

1 We thank all the reviewers for their comments and constructive feedback.

2 **Response to Reviewer 1: Theoretical results:** The reviewer is correct in that some assumptions are needed to derive
3 the results. However, we believe that our results are interesting for two reasons: (i) our paper is the first to present
4 a theoretical analysis of Double Q-learning versus Q-learning that goes beyond asymptotic convergence. (ii) the
5 experimental results (including additional ones we have conducted) support our theory. In particular, as noted in the
6 paper, some of our experiments (e.g., Gridworld) are for models which do not satisfy the assumptions needed for the
7 theory.

8 *Experimental results:* The suggestion to add experiments for the Sutton-Barto example is very interesting. We have
9 indeed performed those simulations now, see Figures (a) and (b). To be consistent with Sutton and Barto, for this
10 example, we have plotted the probability of going left as a function of the number of episodes. Double Q-learning with
11 twice the step-size and averaging does indeed perform better as the theory suggests. We will include the results for
12 these experiments along with the implementation details in the final version of the paper.

13 **Response to Reviewer 2: Theoretical results:** Even though the theory does assume a unique optimal policy, our
14 experimental results indicate that our conclusions hold more broadly. More importantly, we would like to point out that
15 our paper is the first to show that Q-learning will perform strictly better than Double Q-learning (see Theorem 4 in the
16 supplementary material which further strengthens the result in Theorem 2) at least under some nontrivial conditions.
17 Our primary goal is to obtain a theoretical understanding of when and why Double Q-learning performs well. We hope
18 that our idea of analyzing the asymptotic covariance will stimulate further research leading to a deeper understanding of
19 Double Q-learning and its variants.

20 *Experiments beyond linear/tabular settings:* We have now conducted experiments where we have used a neural network
21 to estimate the Q-function. One such result is shown in Figure (b) for the same example suggested by Reviewer 1
22 since it has estimation bias. And again Double Q-learning with twice the step size and averaging performs the best,
23 especially in the initial episodes. We will include the results for these experiments along with the implementation details
24 in the final version of the paper. The comments about Adam/other step-size selection techniques and initialization are
25 interesting, we will explore them in the final version (if the paper is accepted) and/or a longer version of the paper.

26 *References:* We will add and discuss the suggested references in the paper.

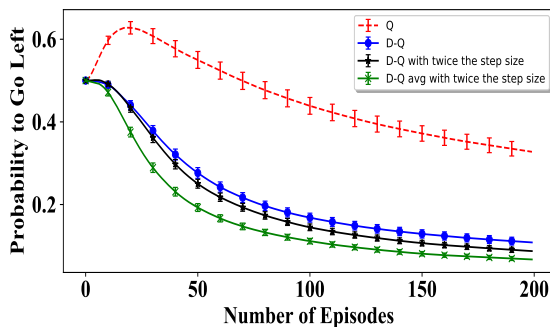
27 **Response to Reviewer 3:**

28 *Experimental validation of the theory:* As mentioned, Figures (a) and (b) present experimental results on an example
29 with estimation bias, where we have also used a neural network for Q-function approximation (Figure (b)). The
30 experimental results in the paper also include examples with non-unique optimal policies. We will include more
31 experimental results in the paper to show robustness with respect to choice of hyper-parameters and also conduct
32 additional experiments in complex RL environments. In the cart pole experiment, we have used the number of episodes
33 to reach a certain reward as our metric and we have the distribution of this quantity. This appears to us to be stronger
34 than knowing just the errors of the means.

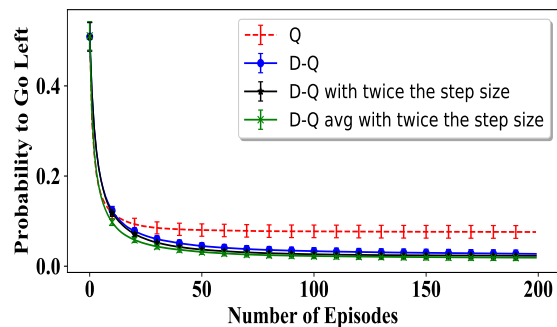
35 *Clarity:* Thanks for the suggestion to include more explanations, we will do so in the final/longer version of the paper.

36 **Response to Reviewer 4:**

37 *Experiments:* Thanks for suggestion about not using code in multiple languages. We will migrate them over to Python
38 before the conference if the paper is accepted. We have now conducted experiments for cases where a neural network is
39 used to approximate the Q-function. Studying such a model theoretically is currently a significant challenge in the field.



(a) Sutton-Barto Tabular



(b) Sutton-Barto Neural Network